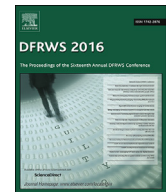


Contents lists available at ScienceDirect

Digital Investigation

journal homepage: www.elsevier.com/locate/diin

DFRWS USA 2016 — Proceedings of the 16th Annual USA Digital Forensics Research Conference

PeekaTorrent: Leveraging P2P hash values for digital forensics



Sebastian Neuner*, Martin Schmiedecker, Edgar R. Weippl

SBA Research, Vienna, Austria

ABSTRACT

Keywords:

Sub-file hashing
Hash-based carving
File whitelisting
p2p file sharing

Sub-file hashing and hash-based carving are increasingly popular methods in digital forensics to detect files on hard drives that are incomplete or have been partially overwritten/modified. While these techniques have been shown to be usable in practice and can be implemented efficiently, they face the problem that a-priori specific “target files” need to be available and at hand. While it is always feasible and, in fact, trivial to create case-specific sub-file hash collections, we propose the creation of case-independent sub-file hash databases. For facilitating hash databases which can be publicly shared among investigators, we propose the usage of data from peer-to-peer file sharing networks such as BitTorrent. Most of the file sharing networks in use today rely on large quantities of hash values for integrity checking and chunk identification, and can be leveraged for digital forensics.

In this paper we show how these hash values can be of use for identifying possibly vast amounts of data and thus present a feasible solution to cope with the ever-increasing case sizes in digital forensics today. While the methodology used is independent of the used file sharing protocol, we harvested information from the BitTorrent network. In total we collected and analyzed more than 3.2 billion hash values from 2.3 million torrent files, and discuss to what extent they can be used to identify otherwise unknown file fragments and data remnants. Using open-source tools like *bulk_extractor* and *hashdb*, these hash values can be directly used to enhance the effectiveness of sub-file hashing at scale.

© 2016 The Author(s). Published by Elsevier Ltd. This is an open access article under the CC BY-NC-ND license (<http://creativecommons.org/licenses/by-nc-nd/4.0/>).

Introduction

One of the current problems in digital forensics is the vast amount of data to be analyzed, as hard drives with 8 terabytes capacity are readily available and the number of devices per person increases steadily. Both are factors for which the current forensic process model does not scale well (Garfinkel, 2010). Acquisition of large data drives can take days, and even though optimization techniques were introduced in the literature recently e.g., sifting collectors

(Richard and Grier, 2015) or file-based deduplication (Neuner et al.), they are not yet used in practice on a larger scale. Slack space, the general availability of counter-forensic tools and increasing importance of RAM content for analysis further challenge the current boundaries of digital forensics. While file whitelisting is a common approach to reduce the number of files to be investigated by an investigator, it is limited in numerous ways: for one, there is currently just one large corpus of hash values which is publicly shared – the NIST National Software Reference Library, containing 43 million file hashes. Secondly, these file hash values rely on hashing an entire file, and are thus unusable for identifying files that are partially modified, or files which have been deleted and partially overwritten.

* Corresponding author.

E-mail addresses: sneuner@sba-research.org (S. Neuner), mschmiedecker@sba-research.org (M. Schmiedecker), eweippl@sba-research.org (E.R. Weippl).

To cope with these problems, we present *PeekaTorrent*, a methodology to identify files and file fragments based on data from publicly available file-sharing networks. It is based on the open-source forensic tools *bulk_extractor* and *hashdb*, and can be readily integrated in forensic processes. It improves the current state-of-the-art on sub-file hashing (Garfinkel and McCarrin, 2015) twofold: for one the hashed sub-file parts are larger than pure sector-based hashes, and thus less prone to false-positives for files that share common data segments. Secondly, we solve the problem that an a-priori sub-file hash database is required by creating one that can be shared openly. Lastly, no participation in file-sharing activity is needed as the torrent metadata or “metainfo”, which is stored in the torrent file, already contains all the necessary information including the sub-file hash values. This information can then be used for file and fragment identification and effective file whitelisting, as well as for other use cases. As such, the contributions of this paper are as follows:

- We present a scalable methodology for identifying files and file fragments based on sub-file hashing and P2P file sharing information.
- We collect and analyze more than 2.3 million torrent files, rendering up to 2.6 petabyte of data identifiable using that information.
- We identify several use cases for file (fragment) identification in the context of both file-whitelisting and blacklisting with that data.
- All obtained data and created source code is available online at <https://www.peakatorrent.org>.

The remainder of this paper is structured as follows: Section 2 provides the necessary background for this paper. Section 3 describes our idea of using sub-file hash values from peer-to-peer file sharing networks in the forensic process, and discusses different use cases where this data can be of value. Section 4 describes our collected data, while the possible benefits are described in Section 5. Section 6 discusses limitations and future work, before we conclude in Section 7.

Background

Digital forensics relies on a multitude of information sources to gain knowledge, ranging from hard drives and file system artefacts (Carrier, 2005), the dynamic content of RAM (Ligh et al., 2014), up to the user files and programs that store information in log files, SQLite databases, or digital images. This leaves the investigator with a broad spectrum of places where to look, and each investigation depends in its specific context and questions to be answered. The general process outline has been defined in both (Brezinski and Killalea, 2002) and (Kent et al.), whereas a great number of current challenges has been discussed in (Garfinkel, 2010). Another problem is the increasing spectrum of devices in use, ranging from smartphones (Hoog, 2011) to smart TVs and numerous other types of devices. Most pressing, however, is the general problem that the average case size is constantly

increasing (Roussev et al., 2013). For one this is due to increasing storage capacities of hard drives, with modern hard drives being able to store many terabytes of data that needs to be analyzed with respect to the traditional approach for digital forensics. Secondly, cloud storage services commonly push information from one device to others automatically, like pictures taken or files edited, leading to duplicate files across devices. Lastly, the density of digital devices surrounding us is increasing, with the average number of devices per user increasing.

Numerous forensic models and publications in recent years were specifically targeted to reduce the manual work needed in investigations with a large amount of data to be analyzed. Among them is the concept of *forensic triage*, which was initially presented in 2006 (Rogers et al., 2006) and more recently quantified regarding the expected amount of computational power needed in (Roussev et al., 2013). The basic idea is that instead of analyzing all the data there is, only a specific subset of files which are known to be of interest are inspected. More recently the concept of *sifting collectors* was proposed (Richard and Grier, 2015) in which the amount of data to be analyzes is reduced by ignoring known areas on hard drives that are of no particular interest, yet still retaining the ability to create bit-identical images if needed. Our approach is different in that it extends the traditional process of forensic imaging by identifying large volumes of both files and file fragments to be either of particular interest (blacklisting), or not of any interest at all as the file is a known-good file (whitelisting).

Both *bulk_extractor* and *hashdb* are two very powerful open-source tools which were published recently by Simon Garfinkel. *Bulk_extractor* (Garfinkel, 2013) recursively scans hard drive content using scanners, and is able to retrieve information in compressed as well as embedded files like PDFs. It is extremely fast, and can use all available cores on a machine to parallelize the task at hand. *Hashdb* (Garfinkel and McCarrin, 2015) uses efficient algorithms to build a lookup database of hash values, much faster than any relational or NoSQL-style database system. It can reliably identify the presence of a given list of target file hash values, and builds on previous work that showed that there is only a small percentage in shared file content on the sector level (Young et al., 2012).

P2P networks for hash values

The basic idea of our approach is to extend the existing knowledge and applicability on sub-file hashing and hash-based carving by leveraging vast amounts of publicly available hash values. While previously hashing was mainly used to uniquely identify entire files of arbitrary size, our concept presented here extends this to hashing variable-sized sub-files portions. Sub-file hashing (Young et al., 2012) as well as hash-based carving (Garfinkel et al., 2010) allow investigators to search for file fragments by hashing either each hard drive sector, or aligned blocks of data. This can also be used if not enough time is available to prove stochastically the presence or absence of specific files e.g. in well below an hour and with only a relatively small error margin. We extend these concepts by mapping sub-file hashes with data from peer-to-peer file sharing

networks with variable block sizes, both usable for black- and white-listing for large volumes of files as well as sampling. We thus extend existing tools and concepts, such as bulk analysis of forensic media using *bulk_extractor* (Garfinkel, 2013) and *hashdb* (Garfinkel and McCarrin, 2015).

Peer-to-peer (P2P) file sharing applications and protocols rely heavily on hashing today, both for integrity and as a foundation for parallelization i.e., downloading multiple parts of a file simultaneously from different users for increased performance. While we used the popular BitTorrent file format for our evaluation, any application that uses sub-file hashing is in many cases directly usable: Dropbox as an example, a popular cloud storage service, hashes blocks of 4 megabytes using SHA-256, and stores them in a local SQLite database (Mulazzani et al., 2011; Kholia and Wegrzyn, 2013). These sub-file hash databases can also be created and maintained privately, say for example based on files and information within a company or an investigative bureau across cases. Our contribution in particular is to propose that these pre-computed hash lists can be used to identify files and sub-files on hard drives. With millions and millions of torrent files publicly shared online, *PeekaTorrent* uses the fact that each and every torrent file indexes all files and also contains their corresponding SHA-1 hash values. For efficiency, the files are split into equally sized pieces, or chunks, solely depending on the overall size of information to be shared (Cohen), in powers of 2 starting with 16 kb. Thus, by splitting the hard drives into equally sized chunks and hashing them using SHA-1, it becomes a matter of comparing hash values to possibly identify hard drive content without relying on file system metadata. Also, this information is freely available without participating in any form of file sharing activities, but leveraging the initial seeders computing power in hashing any form of content.

Torrent files have a rather simple structure (Cohen): they contain generic information like when the torrent was created, which software was used, as well as the specific information of the data to be shared. This includes the size of the blocks, their SHA-1 hash values, and how many of them there are. During the creation of the torrent file, all containing files are concatenated and this stream of data is then split into equally sized blocks of data (except for the last one which does not need to be aligned with the block length). By default, the data is split into 256 kilobyte blocks, but the user can specify arbitrary block sizes during the creation of the torrent file. The size of the torrent file depends mostly on the number of blocks, as for each block it contains a SHA-1 hash value of 20 bytes. To uniquely identify the torrent for both clients and trackers, a SHA-1 hash value is calculated over a subset of the torrents' stored information: the so-called *info_hash*. Fig. 1 shows a graphical representation of the file format, as well as an example from a specific torrent file. The dashed line is the information which is hashed to obtain the *info_hash* value, while for each file the dictionary *files* contains the relative path and the length of the file. *Piece length* is the block size in which the data is split (in the order specified in the *files* field), and the field *pieces* contains the concatenated SHA-1 hash values.

Problem of non-aligned files

One of the problems with the use of torrent files is the way these files are created: prior to hashing all chunks, the files are concatenated (in arbitrary order). For each chunk that contains parts of two files, we cannot use the resulting hash value. This means that only files which are larger than the piece length can be identified, thus biasing the general applicability towards large files (which is obvious when looking at content of file sharing networks). Fig. 2 shows a representation of block hashes in torrents, with the same content as Fig. 1: the SHA-1 value of the first piece is usable, as *codec.exe* spans into the second piece. As such it can be used to uniquely identify that this file has been stored on the hard drive by hashing any hard drive with the same hashing window as the *piece_length* of the torrent. This can be readily integrated into *bulk_extractor*, which already facilitates the necessary requirements by default. If the first file is longer than in our example, and spans e.g., *n* pieces in the torrent file, any of these areas on disc can identify the file as long as the data is stored consecutively somewhere. The second piece in Fig. 2 is not usable for our proposed methodology, as it contains both content from the first and the second file. While it could theoretically happen that the operating system allocates the information in such a way that the hash value could be used, this is not necessarily the case as the files can be stored at different locations on the hard drive and in different ordering. The third piece (the second piece that contains content from *movie.mkv* in our example) is usable if the missing length of the file in the beginning is used for *offset hashing* – it is no longer the *piece_length* which can be used for chunk hashing during acquisition, but rather aligned to the hard drive sectors, which tremendously increases the hash values to be calculated during analysis. Again, this is already integrated in *bulk_extractor*, and the problem remains CPU-bound which means it is solvable if enough computation power is at hand. The hash value for the last piece is unusable, as it must not be of the same length as the others (Cohen) i.e., there is no padding for torrent files.

In the following, we discuss the different use cases where such a vast amount of file fragment information can be of use in the particular context of digital forensics. Other protocols are probably equally suitable, but have not been investigated in detail in this work e.g., Kademlia (Maymounkov and Mazieres, 2002) as well as distributed hash tables in general (Steiner et al., 2009) often use SHA-1 hash values for searching.

Use case 1: file whitelisting

File whitelisting is a well-known technique to identify files during an early phase in digital investigations that are common and of no particular interest. One of the most commonly used database of hash values is the NIST National Software Reference Library (NSRL) reference data set,¹ which comprises at the time of writing of more than 43 million file hash values. Most of these hash values

¹ Online at <http://www.nsl.nist.gov/>.

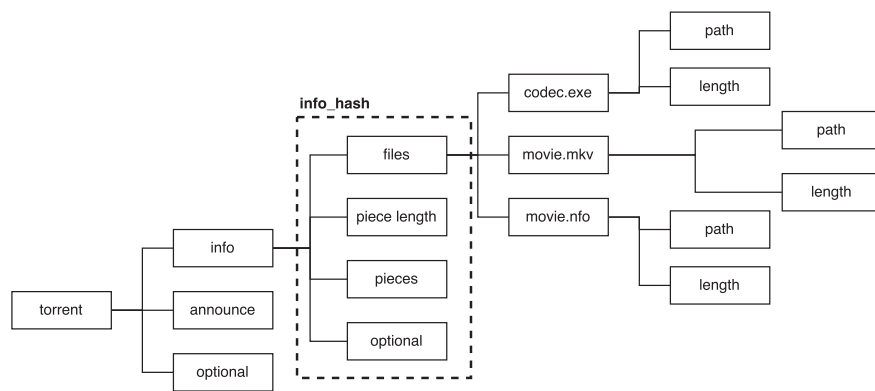


Fig. 1. File content in a torrent file.

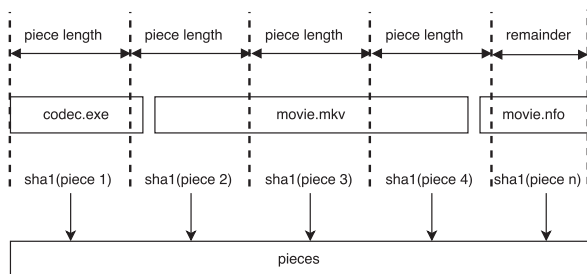


Fig. 2. Chunk hashes.

include binaries and program libraries for software on Windows, whereas our collected data contains information of relevance independent of the used operating system and of much larger file size. While NIST also releases block hash values for the first 4 k and 8 k of about 13 million files, our dataset is able to identify popular files like movies, TV episodes or other commonly shared files on file sharing networks, even if it is deleted and some sectors were already overwritten by the file system.

Use case 2: file blacklisting

File blacklisting is used to find and identify files of particular interest for a specific investigation. While in our evaluation the usability of our data is mostly limited to cases of copyright infringement, it is still of use for investigations in general and might lead to new insights. Nonetheless, building a private sub-file hash database is always a possibility where a script can be used to hash blocks of arbitrary length of e.g. all email attachments in a company, all files on a Sharepoint server or source code within a company. This could also include outright illegal material like pictures and videos related to child pornography. Instead of using perceptual hashing (Breitinger et al) as used by online services like Twitter and Facebook today to detect such files (Ith), sub-file hash values of variable block length can further identify such files without access to such perceptually hashed data.

Use case 3: file fragment identification

File systems in modern operating systems by default, do not overwrite files once they are deleted but rather delete the index pointing to the data or mark the affected storage areas as free-to-use (Carrier, 2005). Depending on the operating system and the file-system in use, as well as the actual user behavior, it is usually not predictable when a specific area will be overwritten. Both methods in our approach described so far work likewise for partially overwritten files, as they do not rely on file system metadata at all. This was already argued in (Young et al., 2012) for sector hashing. As long as the data on disc is not completely overwritten and leaves at minimum the piece length of the torrent files untouched, PeekTorrent will find it.

Shifting the bottleneck

Considering these three use cases, the overall performance scales linearly with the number of available CPU-cores, similar to bulk_extractor. Sub-file hashing can leverage multi-core CPUs and scales with the number of available cores. As the file system metadata is not needed, there is no need for disk seeks. All the data from the hard drive can be split in constant-size chunks, and processed recursively using the hashdb scanner within bulk_extractor.

Evaluation

To evaluate our methodology we implemented numerous steps of the processing outline described above. This includes software we wrote to collect torrent files from the open Internet, and tools to process and use them within the context of a forensic investigation, see <https://www.peakatorrent.org>. This section shows and underlines the applicability of the proposed approach and the methods applied for gathering torrents on a large scale.

Data collection

Collecting a large number of torrents from the open Internet is non-trivial, as new torrents are added constantly and older torrents become unavailable once they are no

longer shared. Only a minority of websites host the torrent files containing all the sub-file hash values themselves, but rather rely on sharing magnet links that point to the information in the completely decentralized *distributed hash tables* (DHTs) (Zhang et al., 2011).

For collecting torrent files we focused on the following three main sources: (i) The Pirate Bay,² (ii) kickassTorrents³ and (iii) various data dumps e.g. from openBay.⁴ For (i) and (ii) we implemented a crawling framework which recursively crawls and parses both websites for every magnet link listed there. After that we extracted the torrent info_hashes from the magnet links, and constructed a download link for the torrent cache website <https://torcache.net/>. For (iii) and those torrent files which weren't hosted at torcache.net we implemented a DHT lookup service, similar to the one Wolchok et al. used in their work (Wolchok and Halderman, 2010). The crawlers for (i) and (ii) were crawling the entire websites, including all sub-categories to get the full archive for a specific point in time (January 2016 in our case).

From the various openBay dumps we were able to extract close to 30 million info_hashes. The dataset from isohunt contained 7.8 million info_hashes, while the complete archive for openBay included 23.5 million hashes. Both data sets were created after the police raid against Pirate Bay in December 2014 caused the website to be shut down. Previously generated data sets also include one notable xml dumps of the Pirate Bay from February 2013 (about 2 million info_hash values). Not all of these files were retrievable using the DHTs, in fact only a small fraction and in particular only newer files. The biggest fraction of torrent files we collected came from kickassTorrents and torcache.net, as torcache.net is used by default to distribute torrent files on behalf of kickassTorrents. So far we have collected 2.3 million torrent files, which we share with the reviewers and later will release them publicly. Our data collection is still going on, and as such the data we collected can be only considered a snapshot in time. Further processing was then done using Python as well as *hashdb*, which was used to efficiently store and query the sub-file hash values.

Theoretic evaluation

Fragmentation of files can be a limiting factor using real cases, as for each time a file is fragmented one chunk (of arbitrary length) is no longer identifiable. Since there is not yet a public instance of a SHA-1 pre-image attack, finding a small number of chunks using PeekaTorrent has a very small likelihood to be coincidentally, and can be used for further analysis steps during the investigation. Compared to previous work (Young et al., 2012; Garfinkel and McCarrin, 2015) the number of false positives is greatly reduced, as the block length used for hashing is larger than the previously used sector/cluster size of 512 or 4096 bytes. Hashing a larger file block e.g., 256 kilobytes, drastically

reduces the probability of resulting in the same hash value (for all files independent of each other). This also implies that shared file content across files, such as the ramping structure for Microsoft Office files as discussed in (Young et al., 2012), is evaded as the block length increases.

Results

Overall, we collected and analyzed more than 2.3 million torrent files. These torrents comprise 3.3 billion block hash values. From these 3.3 billion block hash values, approximately 48% (or 1.62 billion hash block values) are usable to identify million files using various block length. Another 50% (or 1.66 billion hash block values) are usable even though the files do not align with the torrent chunking. 1.1% of the 3.3 billion hash values (or 39 million hash block values) are not usable in our approach, as the blocks and their corresponding hashes comprise content of two or more files. The exact numbers for the most popular torrent block lengths of 2^n (for various n) is shown in Table 1, with exotic chunk sizes omitted ($n = 2871$) for the sake of brevity.

From the 2.3 million torrent files we are able to identify 2.6 petabytes of data using TeekaTorrent, or 32 million files. Regarding only the most common chunk sizes with 100,000 or more torrent files found using our methodology, we are left with 2.1 million torrents. The pre-computed hashdb databases, as well as the raw torrent files and the source code used for this paper can be found on our website <https://www.peekatorrent.org>.

hashdb

We then imported the usable sub-file hash values for all torrents with a piece length of 256 k into hashdb (Garfinkel and McCarrin, 2015). As it can be seen in Table 1, this sums up to 631 million hash values. From these 631 million only 474 million are unique, due to duplicate sub-file hash values. This is due to the fact that the same files can be contained in different torrents, e.g. duplicates for each kickassTorrents and Pirate Bay. Torrent files that became repackaged with different files or file ordering can be another reason to cause this rather large discrepancy. hashdb can then be used to deny that a given sub-file hash value is part of the database using Bloom filters. Otherwise the database is queried, and both filename and info_hash are returned if a corresponding hash value is found. All the features and APIs provided by hashdb are thus fully usable, and the entire project is well documented and active.⁵

While the majority of sub-file hash values are unique within the data we collected (474 million), the long tail of duplicates can be seen in Fig. 3. The x-axis accounts for the number of duplicates found, starting from hash values with 10 duplicates or more. Note that the y-axis is log-scale. In the data there are also 17.8 million distinct sub-file hashes that occurs twice, 2.5 million that occur three times, and about 440,000 that occur four times. We speculate that these hashes are again caused by some form of release

² <https://thepiratebay.se/> and its alternative TLDs.

³ <https://kat.cr/>.

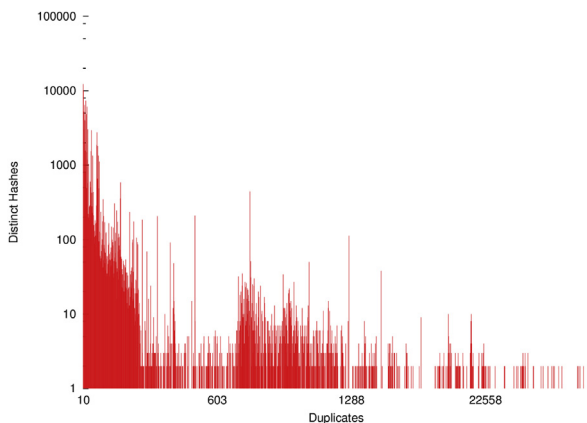
⁴ <https://github.com/isohuntto/openbay-db-dump>.

⁵ <https://github.com/NPS-DEEP/hashdb>.

Table 1

Results of data collection, for 2.3 million torrent files.

Block length	Torrents	Chunks	Usable chunks		Offset chunks		Unusable chunks
16 k	75 k	146 m	123 m	84%	22 m	15%	305 k
32 k	95 k	171 m	112 m	65%	58 m	34%	662 k
64 k	335 k	217 m	124 m	57%	90 m	41%	2 m
128 k	201 k	227 m	115 m	50%	109 m	48%	2 m
256 k	669 k	1.329 b	631 m	47%	690 m	51%	8 m
512 k	297 k	401 m	201 m	50%	194 m	48%	5 m
1024 k	307 k	357 m	165 m	46%	187 m	52%	5 m
2048k	170 k	201 m	75 m	37%	121 m	60%	4 m
4096 k	161 k	229 m	58 m	25%	162 m	70%	8 m
8192 k	18 k	27 m	8 m	30%	17 m	65%	975 k
16384 k	2 k	3 m	315 k	9%	2 m	84%	198 k
Sum:	2.3 m	3.314 b	1.615 b	48%	1.658 b	50%	39 m

**Fig. 3.** Distribution of sub-file hash duplicates.

group information or an embedded URL. The by-far largest number of duplicates observed was caused by one particular hash that occurs 8,462,788 times. We would speculate that this is caused by the “null” hash, for data areas that contain only zeros.

Real runtime on limited hardware

To evaluate our approach furthermore, we took a 5-year old notebook and created a one gigabyte image from a USB thumb drive. The notebook was a Lenovo X200s, with a Core 2 Duo processor (L9400), 4 GB of RAM and a regular hard drive. On the thumb drive we stored the ISO file for the current version of Ubuntu Desktop, which we downloaded over BitTorrent. We created a fresh hashdb database, and seeded it with the extracted SHA-1 hashes of the torrent file. Overall, we extracted 1158 hash values for the Ubuntu image, the chunk size was 512 k. We then used a custom module for bulk_extractor to generate SHA-1 hashes of all blocks bulk_extractor processes, and disabled all other plugins.

Running bulk_extractor with solely the SHA-1 plugin activated on the notebook took 220 s to process the 1 GB image file. Since the CPU has two cores, two threads were spawned to process the image. From the 1158 chunks, 1154 were successfully identified using PeekTorrent. Three

chunks could not be found since the file was stored fragmented in three fragments (verified manually using *fiwalk*), and the last hash value is unusable as it has a different chunk length. Running the same analysis on a modern Xeon with 8 cores plus Hyper-Threading, it took less than 23 s. Running the same image against the hashdb database of all 474 million chunk hashes took 38 s. Since we do not aim to evaluate the performance of either bulk_extractor or hashdb, we do not go into details of further performance numbers. Also, the average fragmentation on hard drives depends heavily on the type of usage, size and operating system. Measuring this for the average case is beyond the scope of this paper.

Discussion

Our results show that a rather large number of block hash values is usable for identifying files based on the data we collected from BitTorrent files, somewhere close to 98%. Due to the nature of file sharing networks and the content distributed there we assume that this is possibly biased due to the fact that these networks commonly share large files like movies in high quality. We did not investigate the distribution of filename, file sizes and to what extent one can expect that the largest file is the first in the torrent file. We assume that this is specific to the application that created the torrent, as this is not specified in the file format of BitTorrent (Cohen).

Half of the usable chunk hashes come with an arbitrary offset due to the placement of the affected files. This is caused by the particularities of BitTorrent files. However, since bulk_extractor processes pages of memory without any file system information, these artefacts are also retrievable (as long as the file is larger than the chunk size). Other sources for sub-file hashing have to be investigated, like other P2P protocols or cloud storage solutions such as Dropbox. We expect similar functionality from other cloud storage solutions like Google Drive, OwnCloud or Microsoft OneDrive as well, where the local data structures could be used as a source for history hash values. Still, using the data we collected we can identify up to 2.6 petabytes of data for 3.3 billion chunks. We expect these values to increase, as we will keep collecting data and publishing it on our website.

Regarding the forensic application and typical use case, many scenarios can come to mind. For one it depends on the data sources used for seeding the sub-file hashing – this can be for example all sent email attachments in a company, a stack of sensitive cooperate documents or encrypted data blobs in the cooperate context. Secondly, this can be easily enlarged by investigators by adding data from private repositories of interesting files, file archives or any other data source at hand like USB thumb drives or portable hard drives, and hashing it in sub-file chunks. Another example could be the cross-linking of files between hard drives: if any of the hard drives during an investigation is hashed with a particular chunk size, all other related drives can be using this information to identify non-fragmented overlaps. After all, this was obviously the original motivation behind the tight connection between *bulk_extractor* and *hashdb*. Foremost, *PeekaTorrent* allows for hard drives without any meta information at all, to find clues on the content as long as the hard drive is not encrypted.

Limitations

While 2.6 petabytes of identifiable files sounds like a lot, its usefulness depends on the particular kind of investigation. If the goal is to whitelist as many files and file fragments as possible on a diverse set of machines, then our approach looks promising. As always in digital forensics, it depends however on the specific context of the investigations, and the questions of interest. For more specific investigations it depends on the type and volume of data – creating sub-file hash values of variable block length is easily scriptable, so if a large repository of files is available, our methodology is applicable. This can be for example all attachments from a mail server, malicious files like malware from anti-virus companies, or even smaller sets of files with a direct connection to an investigation.

Another limitation is the behavior of storage devices, operating systems and file systems: SSDs regularly delete artefacts within the free space using the TRIM command (Bonetti et al., 2013), and depending on the operating system and file system, fragmentation can occur. There are no current numbers on the amount of fragmentation happening, with the latest study on file system metadata being already close to a decade old (Agrawal et al., 2007). Also, the approach only works for files which have at least a file size bigger then the hashing window, or the torrent piece length respectively. Based on our findings with *PeekaTorrent*, only files with a minimal size of 16 kilobytes are identifiable, while the vast amount of files needs to have at least 256 kilobytes due to the nature of the seeding data.

Future work

For future work we plan to evaluate our approach using real hard drives and/or cases. It is generally hard to find representative cases or hard drives, but measuring the applicability of *PeekaTorrent* is our next step. Furthermore, we plan to investigate the usage of GPUs for variable block length hashing. We also plan to make our tools and data

collections more readily applicable, by releasing tools for creating and querying sub-file hash values easy as part of the forensic process. Lastly, our data collection could be enhanced by focusing on popular file torrents, and by collecting more files over time (which is expected to continue for the near future) and from additional torrent websites as well as from DHT crawlers.

Conclusion

In this paper we have demonstrated how vast amounts of sub-file hash values can be of use in digital forensics. We evaluated the idea using torrent files from popular file sharing platforms, and collected more than 2.3 million torrent files for our analysis. Based on these torrent files we extracted more then 3 billion SHA-1 sub-file hash values, and are able to identify up to 32 million files or 2.6 petabytes of information using this data set. Both, the collected data as well as the written software tools are available under open source licenses.

Acknowledgements

We thank our shepherd Judson Powers for guiding us to a highly improved version of the initial paper. We owe particular thanks to our student Daniel Gasperschitz, for writing the SHA-1 module for *bulk_extractor*. This research was supported by the Austrian Research Promotion Agency (FFG) through the Bridge Early Stage grant P846070 (SpeedFor) and the COMET K1 program.

References

- Agrawal N, Bolosky WJ, Douceur JR, Lorch JR. A five-year study of file-system metadata. *ACM Transactions on Storage (TOS)* 2007;3(3):9.
- Bonetti G, Viglione M, Frossi A, Maggi F, Zanero S. A comprehensive black-box methodology for testing the forensic characteristics of solid-state drives. In: *Proceedings of the 29th annual computer security applications conference*, ACM; 2013. p. 269–78.
- F. Breiting, B. Guttman, M. McCarrin, V. Roussev, D. White, Approximate matching: definition and terminology, URL: <http://csrc.nist.gov/publications/drafts/800-168/sp800-168-draft.pdf>.
- Brezinski D, Killalea T. Guidelines for evidence collection and archiving. RFC 3227 (Best Current Practice). Feb. 2002.
- Carrier B. *File system forensic analysis*. Addison-Wesley Professional; 2005.
- Garfinkel SL. Digital forensics research: the next 10 years. *Digit Investig* 2010;7:S64–73.
- Garfinkel SL. Digital media triage with bulk data analysis and *bulk_extractor*. *Comput Secur* 2013;32:56–72.
- Garfinkel SL, McCarrin M. Hash-based carving: searching media for complete files and file fragments with sector hashing and *hashdb*. *Digit Investig* 2015;14:S95–105.
- Garfinkel S, Nelson A, White D, Roussev V. Using purpose-built functions and block hashes to enable small block and sub-file forensics. *Digit Investig* 2010;7:S13–23.
- Hoog A. *Android forensics: investigation, analysis and mobile security for Google Android*. Elsevier; 2011.
- T. Ith, Microsoft's photodna: protecting children and businesses in the cloud, online at: <https://news.microsoft.com/features/microsofts-photodna-protecting-children-and-businesses-in-the-cloud/> [15.07.15].
- K. Kent, T. Grance, H. Dang, Nist special publication 800-86, guide to integrating forensic techniques into incident response.
- Kholia D, Wegrzyn P. Looking inside the (drop) box. In: *7th USENIX workshop on offensive technologies (WOOT)*; 2013.
- Ligh MH, Case A, Levy J, Walters A. *The art of memory forensics: detecting malware and threats in windows, linux, and Mac memory*. John Wiley & Sons; 2014.

- Maymounkov P, Mazières D. Kademlia: a peer-to-peer information system based on the xor metric. In: *Peer-to-Peer systems*. Springer; 2002. p. 53–65.
- Mulazzani M, Schrittwieser S, Leithner M, Huber M, Weippl ER. Dark clouds on the horizon: using cloud storage as attack vector and online slack space. In: *USENIX security symposium*, San Francisco, CA, USA; 2011. p. 65–76.
- S. Neuner, M. Schmiedecker, E. Weippl, Effectiveness of file-based deduplication in digital forensics, *Security and Communication Networks*.
- B. Cohen, The bittorrent protocol specification, bep-3, online at: http://www.bittorrent.org/beps/bep_0003.html.
- Richard III GG, Grier J. Rapid forensic acquisition of large media with sifting collectors. *Digit Investig* 2015;14:S34–44.
- Rogers MK, Goldman J, Mislan R, Wedge T, Debroya S. Computer forensics field triage process model. In: *Proceedings of the conference on digital forensics, security and law*, association of digital forensics, security and law; 2006. p. 27.
- Roussev V, Quates C, Martell R. Real-time digital forensics and triage. *Digit Investig* 2013;10(2):158–67.
- Steiner M, En-Najjary T, Biersack EW. Long term study of peer behavior in the kad dht. *IEEE/ACM Trans Netw (TON)* 2009;17(5):1371–84.
- Wolchok S, Halderman JA. Crawling bittorrent dhts for fun and profit. In: *4th USENIX workshop on offensive technologies (WOOT)*; 2010.
- Young J, Foster K, Garfinkel S, Fairbanks K. Distinct sector hashes for target file detection. *Computer* 2012;12:28–35.
- Zhang C, Dhungel P, Wu D, Ross KW. Unraveling the bittorrent ecosystem. *Parallel Distributed Syst IEEE Trans* 2011;22(7):1164–77.